



# Learning to Pre-train Graph Neural Networks

Yuanfu Lu<sup>1</sup>, Xunqiang Jiang<sup>1</sup>, Yuan Fang<sup>2</sup>, Chuan Shi<sup>1</sup>

<sup>1</sup>Beijing University of Posts and Telecommunications

<sup>2</sup>Singapore Management University



北京郵電大學

Beijing University of Posts and Telecommunications



SMU

SINGAPORE MANAGEMENT  
UNIVERSITY

## GNNs

- ▶ node-level representation

$$\begin{aligned}\mathbf{h}_v^l &= \Psi(\psi; \mathcal{A}, \mathcal{X}, \mathcal{Z})^l \\ &= \text{UPDATE}(\mathbf{h}_v^{l-1}, \\ &\quad \text{AGGREGATE}(\{(\mathbf{h}_v^{l-1}, \mathbf{h}_u^{l-1}, \mathbf{z}_{uv}) : u \in \mathcal{N}_v\}))\end{aligned}$$

- ▶ graph-level representation

$$\mathbf{h}_G = \Omega(\omega; \mathbf{H}^l) = \text{READOUT}(\{\mathbf{h}_v^l | v \in \mathcal{V}\})$$

## Pre-train GNNs

- ▶

$\theta_0$  is pre-trained

without accommodating the adaptation in fine-tuning

$$\theta_0 = \arg \min_{\theta} \mathcal{L}^{pre}(f_{\theta}; \mathcal{D}^{pre})$$

$$\theta_1 = \theta_0 - \eta \nabla_{\theta_0} \mathcal{L}^{fine}(f_{\theta_0}; \mathcal{D}^{tr})$$

*learn how to pre-train GNNs*

▶ How to **narrow the gap** caused by different optimization objectives?

- ▶ *SOTAs fall into a two-step paradigm with a gap*
- ▶ *Solution: **learn to pre-train** (meta learning)*

▶ How to **simultaneously preserve node- and graph-level information?**

- ▶ *SOTAs either only consider the node-level pre-training or require supervised information for graph-level pre-training*
- ▶ *Solution: **intrinsic self-supervision***

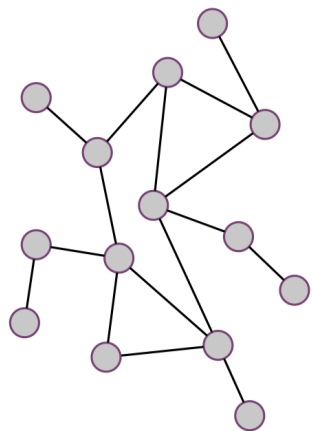
Pre-train a GNN model over a graph  $\mathcal{G} \in \mathcal{D}^{pre}$

- ▶ sample sub-structures  $\mathcal{D}_{\mathcal{G}}^{tr}$  for training  
(the training data of a *simulated downstream task*)
- ▶ **mimic the evaluation** on testing sub-structures  $\mathcal{D}_{\mathcal{G}}^{te}$

$$\theta_0 = \arg \min_{\theta} \sum_{\mathcal{G} \in \mathcal{D}^{pre}} \mathcal{L}^{pre}(f_{\theta - \alpha \nabla_{\theta} \mathcal{L}^{pre}(f_{\theta}; \mathcal{D}_{\mathcal{G}}^{tr})}; \mathcal{D}_{\mathcal{G}}^{te})$$

the fine-tuned parameters

(in a similar manner as the fine-tuning step on the downstream task)

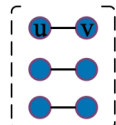


$$\mathcal{G} = \{\mathcal{V}, \mathcal{E}, \mathcal{X}, \mathcal{Z}\}$$

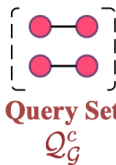
(a) An Example of Graph

**Parent Task**  
 $\mathcal{T}_{\mathcal{G}} = \{\mathcal{T}_{\mathcal{G}}^1, \dots, \mathcal{T}_{\mathcal{G}}^k\}$

**Child Task**  $\mathcal{T}_{\mathcal{G}}^c$

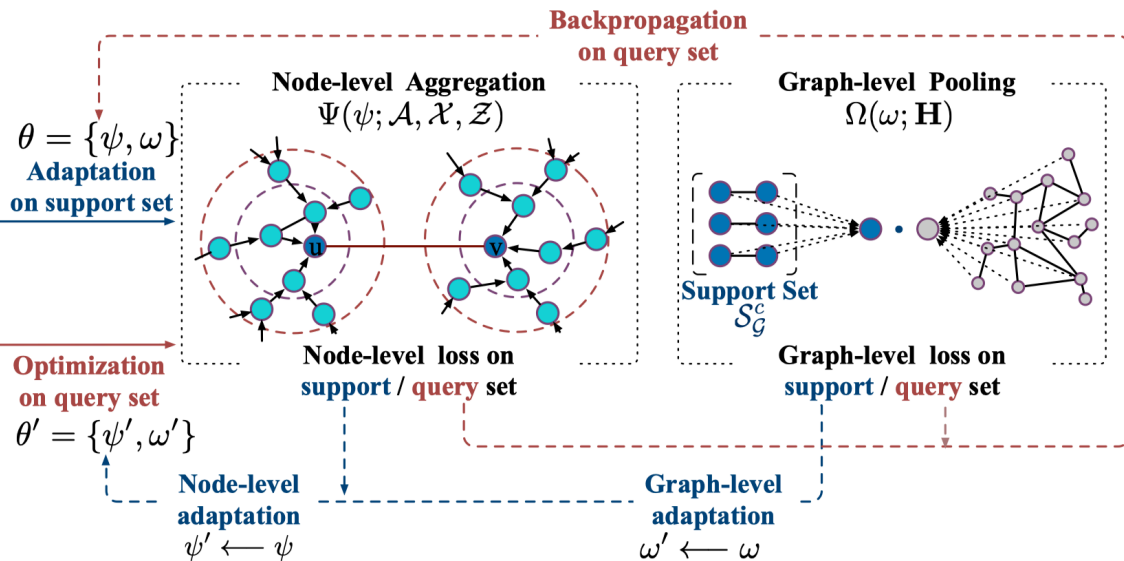


**Support Set**  
 $S_{\mathcal{G}}^c$



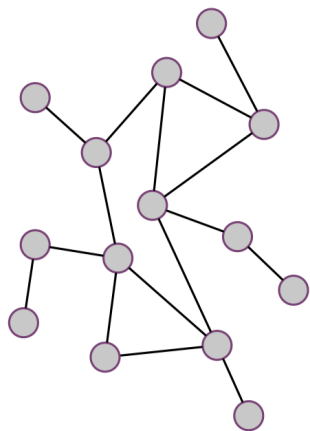
**Query Set**  
 $Q_{\mathcal{G}}^c$

(b) Task Construction



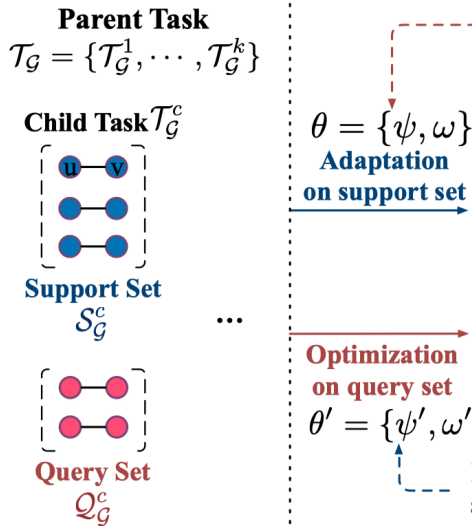
(c) Dual Adaptation in Self-supervised Base Model

## Task Construction



$$\mathcal{G} = \{\mathcal{V}, \mathcal{E}, \mathcal{X}, \mathcal{Z}\}$$

(a) An Example of Graph



(b) Task Construction

- ▶ the pre-training data  
 $\mathcal{D}^{pre} = \{\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_N\}$

- ▶ A task involving a graph

$$\mathcal{T}_{\mathcal{G}} = (\mathcal{S}_{\mathcal{G}}, \mathcal{Q}_{\mathcal{G}})$$

- ▶ *gradient descent* w.r.t. the loss on  $\mathcal{S}_{\mathcal{G}}$
- ▶ *optimize* the performance on  $\mathcal{Q}_{\mathcal{G}}$
- ▶ *simulating the training and testing* in the fine-tuning step

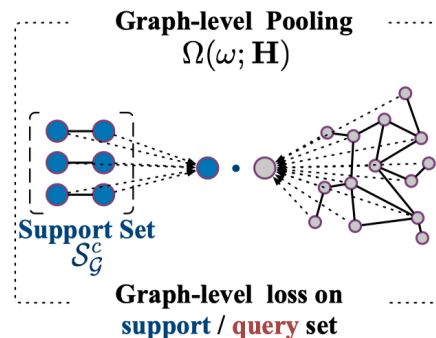
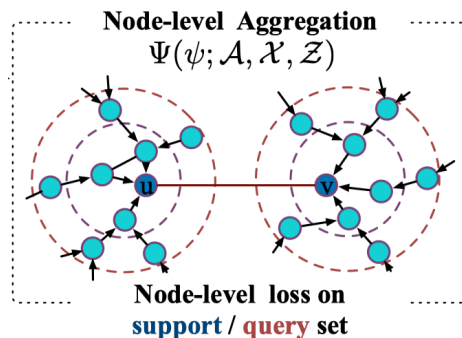
## Self-supervised Base Model

### ► node-level aggregation

$$\mathcal{L}^{node}(\psi; \mathcal{S}_G^c) = \sum_{(u,v) \in \mathcal{S}_G^c} -\ln(\sigma(\mathbf{h}_u^\top \mathbf{h}_v)) - \ln(\sigma(-\mathbf{h}_u^\top \mathbf{h}_{v'}))$$

### ► graph-level pooling

$$\mathcal{L}^{graph}(\omega; \mathcal{S}_G) = \sum_{c=1}^k -\log(\sigma(\mathbf{h}_{\mathcal{S}_G^c}^\top \mathbf{h}_G)) - \log(\sigma(-\mathbf{h}_{\mathcal{S}_G^c}^\top \mathbf{h}_{G'}))$$



$$\mathcal{L}_{\mathcal{T}_G}(\theta; \mathcal{S}_G) = \mathcal{L}^{graph}(\omega; \mathcal{S}_G) + \frac{1}{k} \sum_{c=1}^k \mathcal{L}^{node}(\psi; \mathcal{S}_G^c)$$



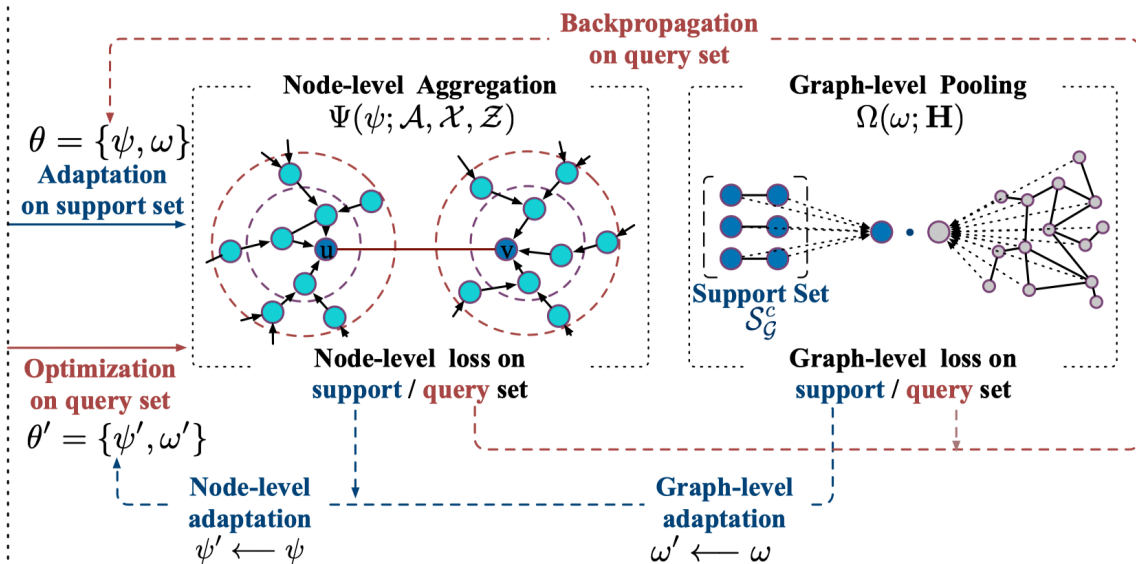
## Dual Adaptation

### ► node-level adaptation

$$\psi' = \psi - \alpha \frac{\partial \sum_{c=1}^k \mathcal{L}^{node}(\psi; \mathcal{S}_G^c)}{\partial \psi}$$

### ► graph-level adaptation

$$\omega' = \omega - \beta \frac{\partial \mathcal{L}^{graph}(\omega; \mathcal{S}_G)}{\partial \omega}$$



(c) Dual Adaptation in Self-supervised Base Model

$$\theta \leftarrow \theta - \gamma \frac{\partial \sum_{G \in \mathcal{D}^{pre}} \mathcal{L}_{T_G}(\theta'; \mathcal{Q}_G)}{\partial \theta}$$

## A new dataset for pre-training GNNs

### ► Datasets

Dataset	Biology	PreDBLP
#subgraphs	394,925	1,054,309
#labels	40	6
#subgraphs for pre-training	306,925	794,862
#subgraphs for fine-tuning	88,000	299,447

### ► Baselines

- EdgePred *to predict the connectivity of node pairs*
- DGI *to maximize mutual information across the graph's patch representations*
- ContextPred *to explore graph structures*
- AttrMasking *to learn the regularities of node/edge attributes*

### ► GNN Architectures

- GCN, GraphSAGE, GAT, GIN

# Performance Comparison



北京邮电大学  
Beijing University of Posts and Telecommunications



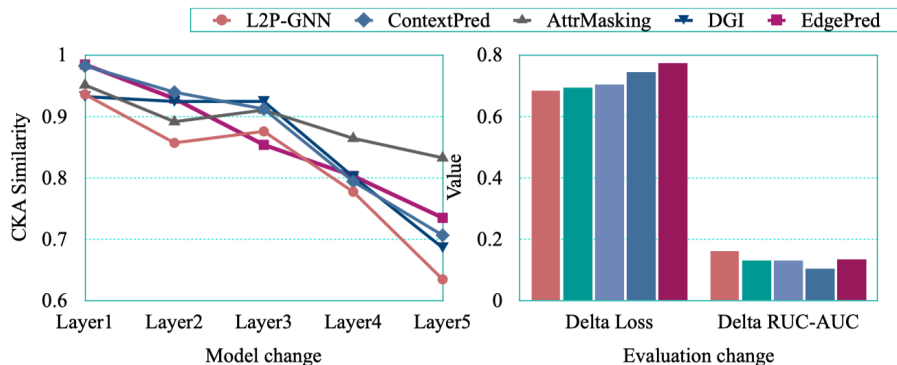
SMU  
SINGAPORE MANAGEMENT  
UNIVERSITY



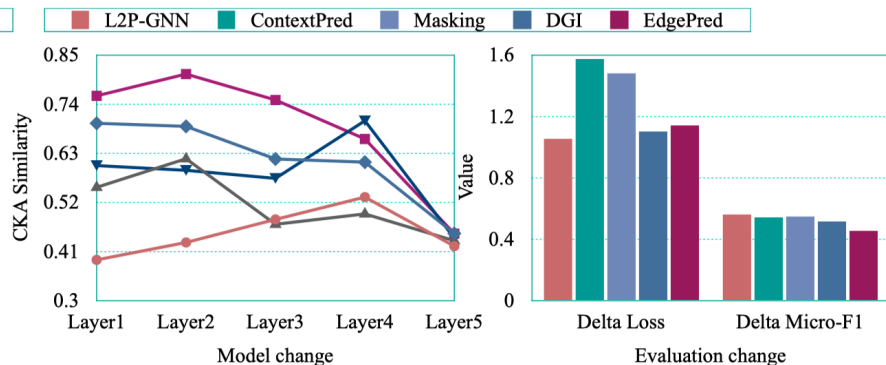
Table 2: Experimental results (mean  $\pm$  std in percent) of different pre-training strategies w.r.t. various GNN architectures. The improvements are relative to the respective GNN without pre-training.

Model	Biology				PreDBLP			
	GCN	GraphSAGE	GAT	GIN	GCN	GraphSAGE	GAT	GIN
No pre-train	63.22 $\pm$ 1.06	65.72 $\pm$ 1.23	68.21 $\pm$ 1.26	64.82 $\pm$ 1.21	62.18 $\pm$ 0.43	61.03 $\pm$ 0.65	59.63 $\pm$ 2.32	69.01 $\pm$ 0.23
EdgePred	64.72 $\pm$ 1.06	67.39 $\pm$ 1.54	67.37 $\pm$ 1.31	65.93 $\pm$ 1.65	65.44 $\pm$ 0.42	63.60 $\pm$ 0.21	55.56 $\pm$ 1.67	69.43 $\pm$ 0.07
DGI	64.33 $\pm$ 1.14	66.69 $\pm$ 0.88	68.37 $\pm$ 0.54	65.16 $\pm$ 1.24	65.57 $\pm$ 0.36	63.34 $\pm$ 0.73	61.30 $\pm$ 2.17	69.34 $\pm$ 0.09
ContextPred	64.56 $\pm$ 1.36	66.31 $\pm$ 0.94	66.89 $\pm$ 1.98	65.99 $\pm$ 1.22	66.11 $\pm$ 0.16	62.55 $\pm$ 0.11	58.44 $\pm$ 1.18	69.37 $\pm$ 0.21
AttrMasking	64.35 $\pm$ 1.23	64.32 $\pm$ 0.78	67.72 $\pm$ 1.16	65.72 $\pm$ 1.31	65.49 $\pm$ 0.52	62.35 $\pm$ 0.58	53.34 $\pm$ 4.77	68.61 $\pm$ 0.16
L2P-GNN (Improv.)	<b>66.48</b> $\pm$ 1.59 (5.16%)	<b>69.89</b> $\pm$ 1.63 (6.35%)	<b>69.15</b> $\pm$ 1.86 (1.38%)	<b>70.13</b> $\pm$ 0.95 (8.19%)	<b>66.58</b> $\pm$ 0.28 (7.08%)	<b>65.84</b> $\pm$ 0.37 (7.88%)	<b>62.24</b> $\pm$ 1.89 (4.38%)	<b>70.79</b> $\pm$ 0.17 (2.58%)

- ▶ **6.27% and 3.52%** improvements compared to the best baseline
- ▶ **8.19% and 7.88%** gains relative to non-pretrained models
- ▶ **negative transfer** harms the generalization of the pre-trained GNNs (e.g., EdgePred and AttrMasking strategies w.r.t. GAT)



(a) Biology dataset

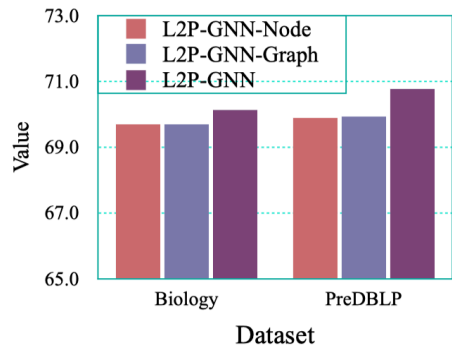


(b) PreDBLP dataset

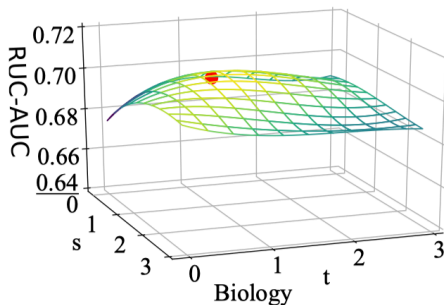
## Comparative Analysis

*whether L2P-GNN narrows the gap between pre-training and fine-tuning?*

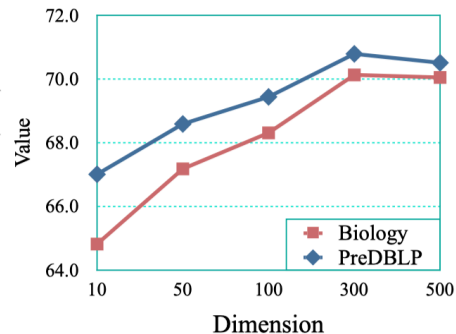
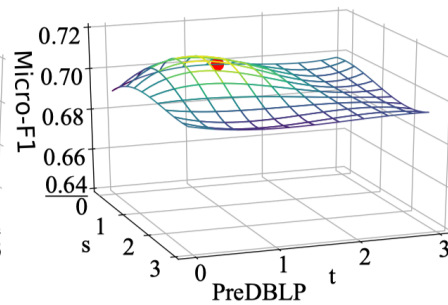
- ▶ Comparison of the pre-trained GNN model **before and after fine-tuning**
  - ▶ **Centered Kernel Alignment (CKA) similarity** between the parameters
    - ▶ *Smaller similarity, larger changes of model parameters*
  - ▶ changes in loss and performance (**delta loss and RUC-AUC/Micro-F1**)
    - ▶ *Smaller change, more easily achieve the optimal point*



(a) Ablation study.



(b) Node- and graph-level adaptation steps ( $s, t$ ).



(c) Dimension analysis.

## ► Ablation Study

- L2P-GNN-Node with **only node-level** adaptation
- L2P-GNN-Graph with **only graph-level** adaptation

## ► Parameter Analysis

- the number of node- and graph-level **adaptation steps** ( $s, t$ )
- the **dimension** of node representations

- ▶ **A problem**

- ▶ *there exists a divergence between the pre-training and fine-tuning objectives, resulting in suboptimal pre-trained GNN models*

- ▶ **A solution**

- ▶ *a self-supervised pretraining strategy for GNNs, L2P-GNN, which attempts to learn how to fine-tune during the pre-training process in the form of transferable prior knowledge*

- ▶ **A dataset**

- ▶ *a new large-scale graph structured data for pre-training GNNs*

# Thank you !

## Q&A

Codes and datasets: <https://github.com/rootlu/L2P-GNN>



北京郵電大學

Beijing University of Posts and Telecommunications



SMU

SINGAPORE MANAGEMENT  
UNIVERSITY